

AI Fundamentals of Counterfactuals

Computational Foundations of Counterfactuals



Plan for the Day

CF Tutorial	
TIME	Topics
9:00 AM	Introduction
	<i>Hello and Introducing Ourselves!</i>
	Hands-on: <i>Trying Our Study (follow link)</i>
9:30 AM	Historical Fundamentals of Counterfactuals
	<i>From Philosophy to XAI (via Psychology)</i>
	<i>Two Sample User Studies and Q&A</i>
10:30 AM	COFFEE (10:30-11:00)
11:00 AM	Fundamentals of Counterfactuals in AI
	<i>Formalisation</i>
	<i>Modelling Approaches & Key Constraints</i>
11:30 AM	Using Counterfactual Algorithms
	Hands-on: <i>A Counterfactual Toolbox (AA)</i>
	Hands-on: <i>Checking Out Notebooks and Q&A</i>
12:00 PM	Fundamentals of User Studies
	<i>User Studies I: A Simple Two-Group Design</i>
12:30 PM	LUNCH (12:30-14:00)
2:00 PM	Algorithmic Growth Points
	<i>Computational Future Directions and Q&A</i>
2:30 PM	More Fundamentals of User Studies
	<i>User Studies I: A Simple Two-Group Design (cont.)</i>
3:00 PM	COFFEE (15:00-15:30)
3:30 PM	From Fundamentals to an Actual User Study
	<i>User Studies II: A More Complex Design</i>
	<i>User Studies III: Even More Complex Designs</i>
	Hands-on: <i>Looking At Our Study</i>
5:00 PM	Closing Session, Discussion and Final Q&A
	TUTORIAL END



You Are Here!

Detailed Outline



Formalization



Modeling Approaches



Key Constraints

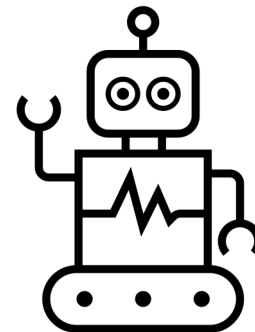


Hands-On: Counterfactuals in Python

Formalization of Counterfactuals

- Two properties:
 1. Contrastive
 2. Proximity/Closeness/Cost
 - Often a p-norm – but is this realistic?
=> Domain specific!

Y, If X had happened!



Modeling Approaches

- Optimization problem [\[Wachter 2017\]](#):

$$\arg \min_{\vec{x}_{cf} \in \mathbb{R}^d} \boxed{\ell(h(\vec{x}_{cf}), y_{cf})} + C \cdot \boxed{\theta(\vec{x}_{cf}, \vec{x}_{orig})}$$

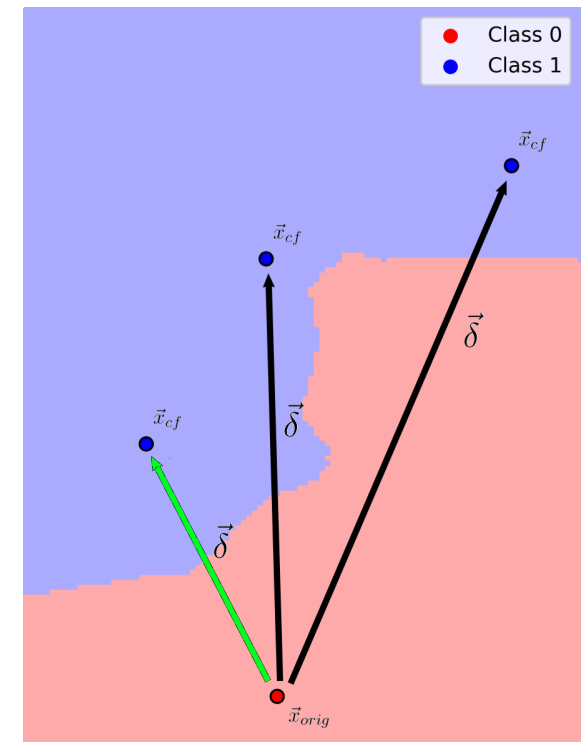
↑
Contrastive
↑
Cost/Proximity

Change $\vec{\delta} = \vec{x}_{cf} - \vec{x}_{orig}$

In constraint form:

$$\arg \min_{\vec{x}_{cf} \in \mathbb{R}^d} \boxed{\theta(\vec{x}_{cf}, \vec{x}_{orig})} \quad \text{s.t.} \quad \boxed{h(\vec{x}_{cf}) = y_{cf}}$$

↑
Cost/Proximity
↑
Contrastive

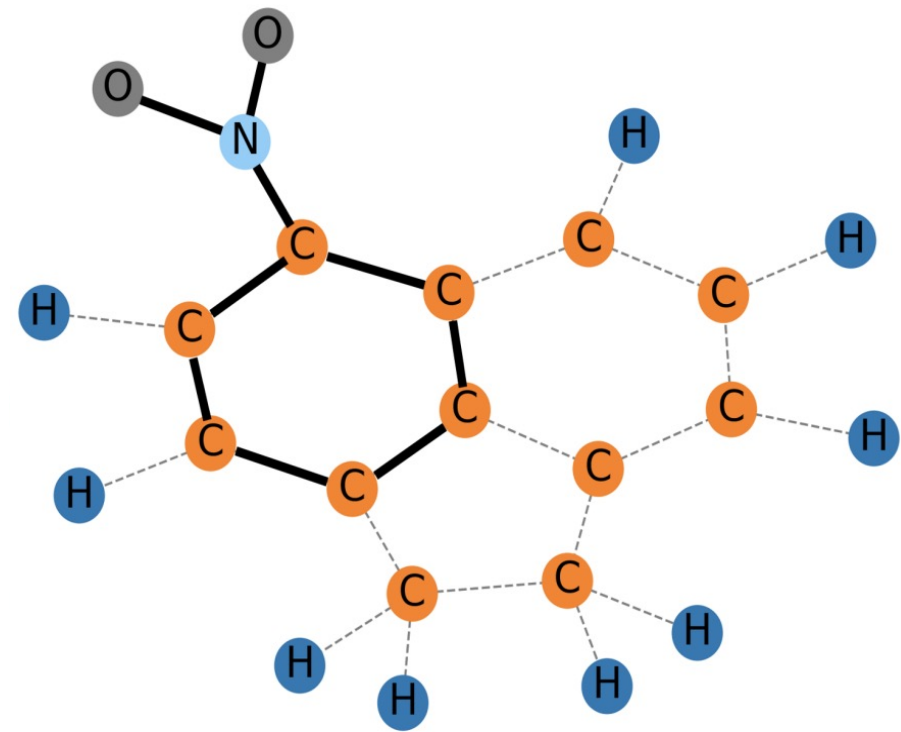


Modeling Approaches

- How to solve the optimization problems?

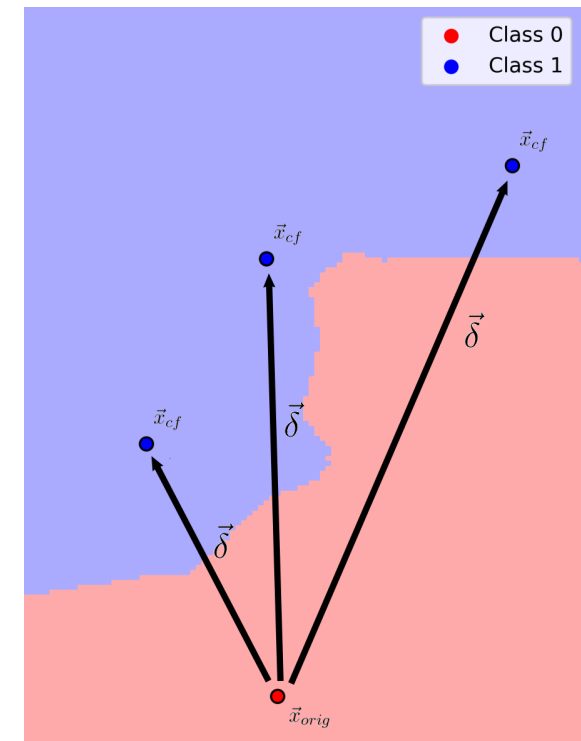
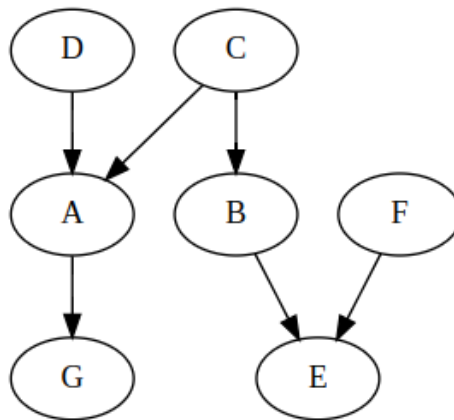
$$\arg \min_{\vec{x}_{cf} \in \mathbb{R}^d} \ell(h(\vec{x}_{cf}), y_{cf}) + C \cdot \theta(\vec{x}_{cf}, \vec{x}_{orig})$$

- Black-box solvers
- Gradient descent
- Domain-specific methods (images, text, graphs, etc.)



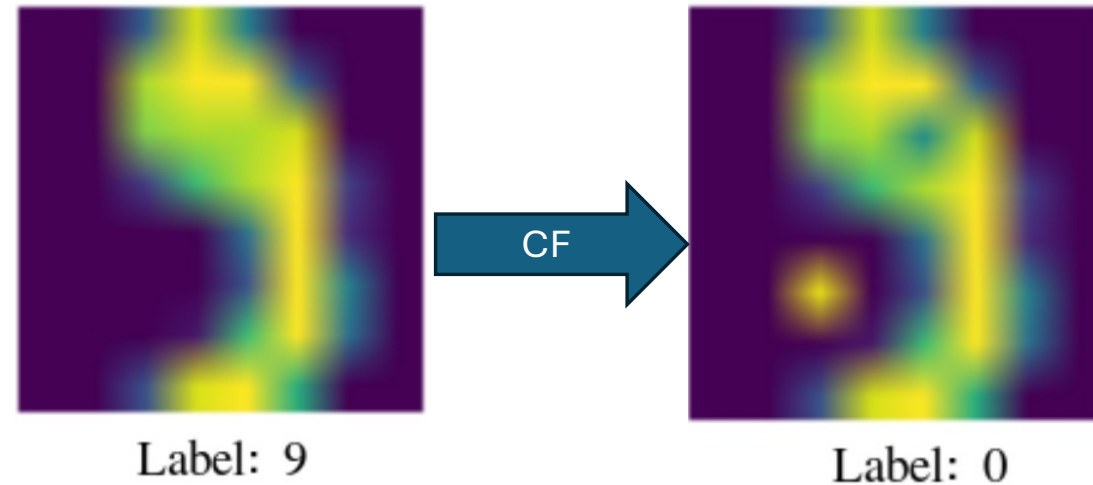
Things to keep in Mind

- Uniqueness?
=> "Rashomon Effect"
- What about Causality? [\[Karimi 2021\]](#)
=> Often, feature independence is assumed!
Incorporate Causal models!



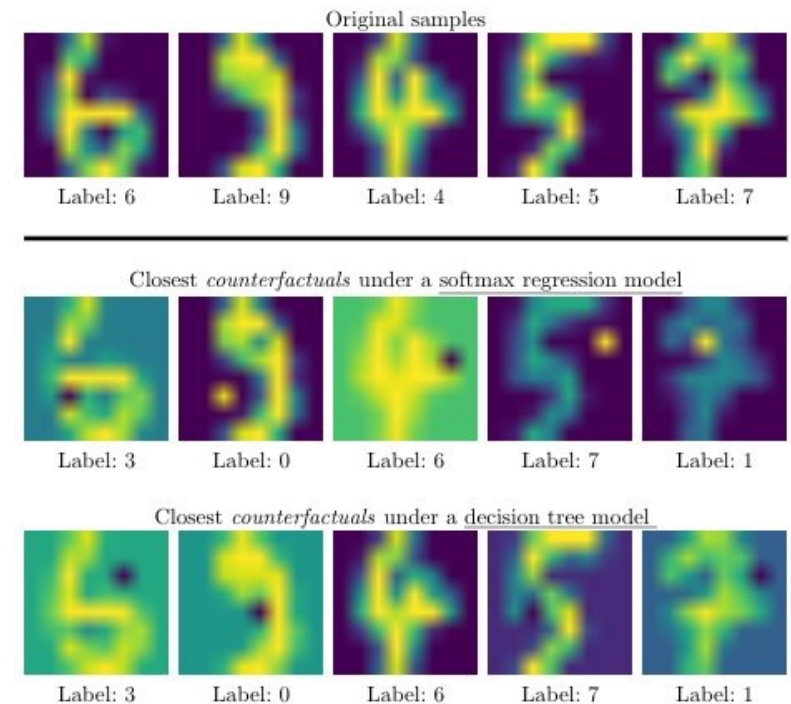
Fundamentals & Key Constraints

- Contrastivity and Cost/Proximity are not enough!
 - e.g. CF \neq Adversarial
- Other important aspects:
 - Plausibility
 - Actionability
 - Diversity
 - ...



Plausibility & Actionability

- Problem: Classic CFs [\[Wachter 2017\]](#) often adversarial
=> Not useful in practice! [\[Smyth 2022\]](#)
- **CFs must be plausible and actionable**

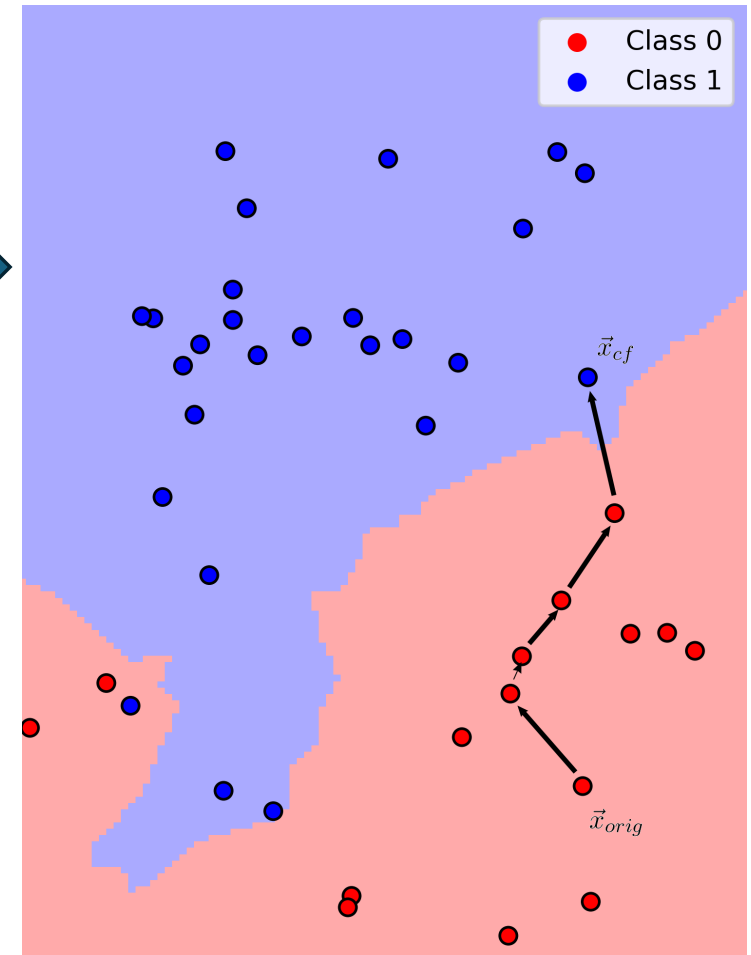


Plausibility & Actionability

- Two popular & general approaches:

- FACE [\[Poyiadzi 2020\]](#)
- Density constraints [\[Artelt 2020\]](#)

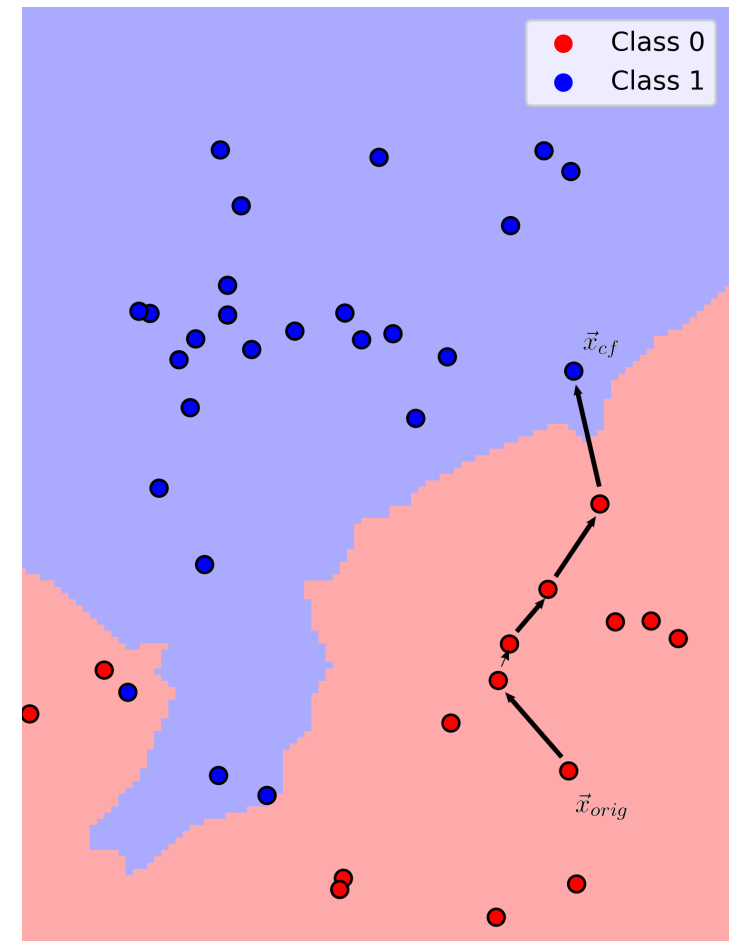
$$\begin{aligned} & \arg \min_{\vec{x}_{cf} \in \mathbb{R}^d} \theta(\vec{x}_{cf}, \vec{x}_{orig}) \\ & \text{s.t. } h(\vec{x}_{cf}) = y_{cf} \\ & \hat{p}_{y_{cf}}(\vec{x}_{cf}) \geq \underbrace{\delta}_{\text{lower-bound on density}} \end{aligned}$$



FACE: Feasible, Actionable CFs

[Poyiadzi 2020]

- Sequence of changes
- Limit feasible set to observed data
 - "No feature independence"
 - Actionable & Plausible
- Graph: KDE, k-NN, or ϵ -graph
 - => Find shortest path

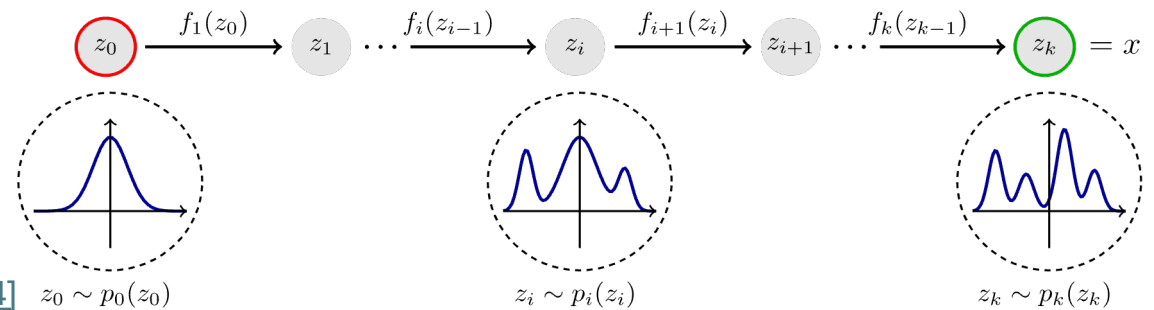


Density Constraints

- **CF in high-density region**

- Methods:

- GMMs [\[Artelt 2020\]](#)
- KDE [\[Förster 2021\]](#)
- GANs [\[Mertes 2022\]](#)
- Normalizing Flows [\[Wielopolski 2024\]](#)



$$\arg \min_{\vec{x}_{\text{cf}} \in \mathbb{R}^d} \theta(\vec{x}_{\text{cf}}, \vec{x}_{\text{orig}})$$

$$\text{s.t. } h(\vec{x}_{\text{cf}}) = y_{\text{cf}}$$

$$\hat{p}_{y_{\text{cf}}}(\vec{x}_{\text{cf}}) \geq$$

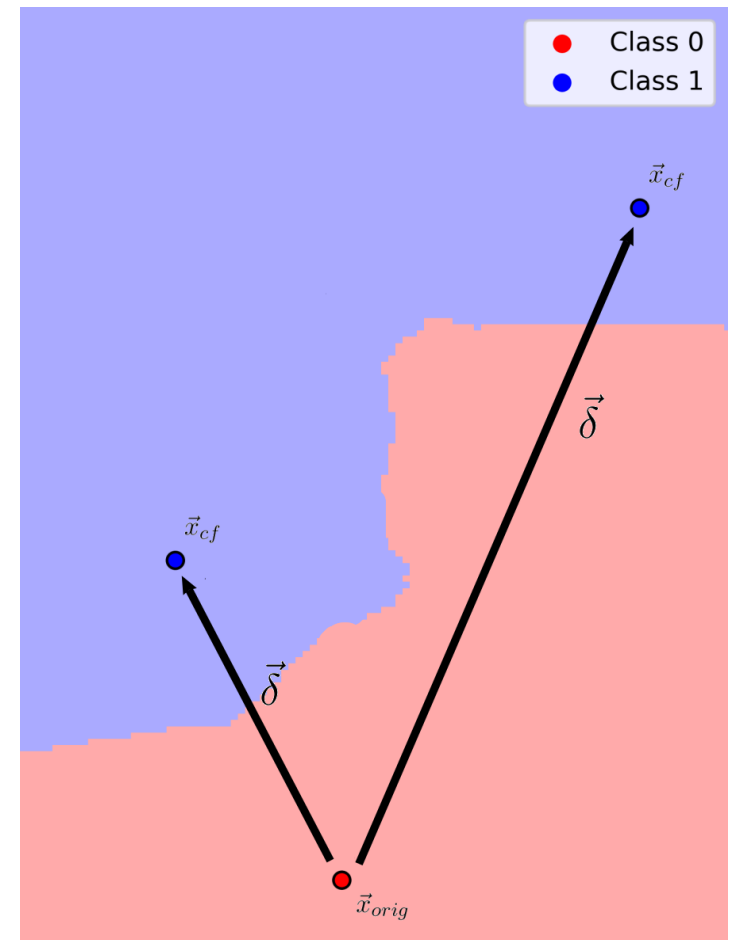
$$\underbrace{\delta}_{\text{lower-bound on density}}$$

lower-bound on density

Diversity

- Missing uniqueness
=> "Rashomon Effect"
- Generate multiple "different" CFs
 - Let the user choose!
 - More robust? [\[Leofante 2024\]](#)

RQ: How to model diversity?



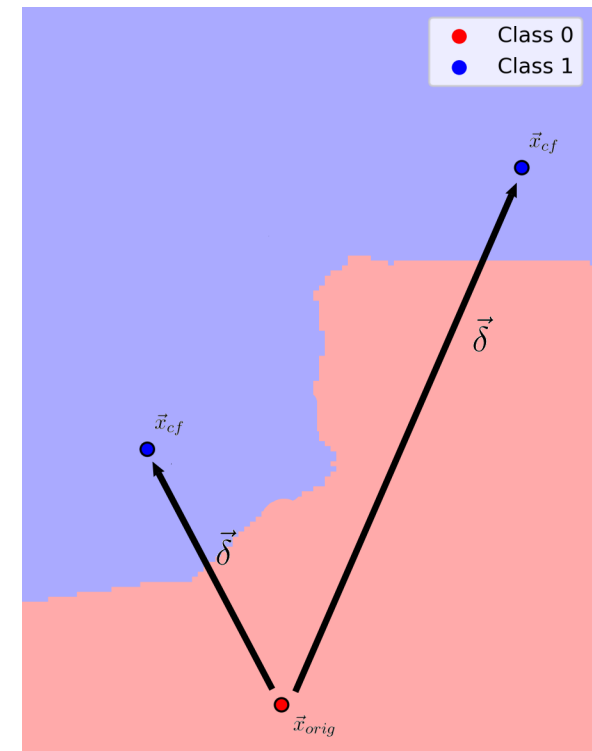
Diversity

- Yet another constraint:
 - Quantify difference between CFs (e.g. cost, common features, etc.)
- Toolbox: DiCE [\[Mothilal 2020\]](#)
=> See Hands-on session

$$\arg \min_{\{\vec{x}_{cf_i} \in \mathbb{R}^d\}} \underbrace{\sum_i \ell(h(\vec{x}_{cf_i}), y_{cf})}_{\text{Constrasting}} + \underbrace{C_1 \sum_i \theta(\vec{x}_{cf_i}, \vec{x}_{orig})}_{\text{Complexity}} + \underbrace{C_2 \sum_{i,j} \psi(\vec{x}_{cf_i}, \vec{x}_{cf_j})}_{\text{Diversity}}$$

Summary & Conclusion

- Two important properties:
 - Contrastive
 - Cost/Proximity
- Key constraints:
 - Plausibility & Actionability
 - Diversity



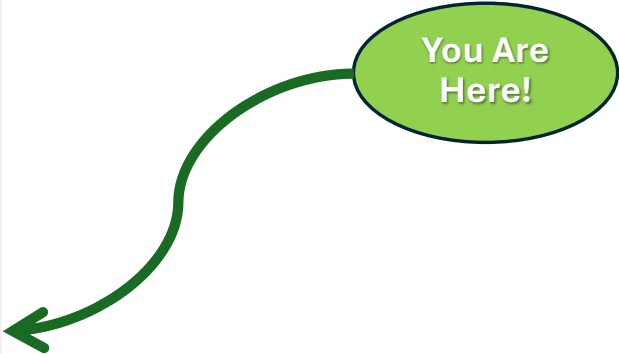
Hands-On: Counterfactuals in Python



<https://tinyurl.com/6mc35pya>

Plan for the Day

CF Tutorial	
TIME	Topics
9:00 AM	Introduction
	<i>Hello and Introducing Ourselves!</i>
	Hands-on: <i>Trying Our Study (follow link)</i>
9:30 AM	Historical Fundamentals of Counterfactuals
	<i>From Philosophy to XAI (via Psychology)</i>
	<i>Two Sample User Studies and Q&A</i>
10:30 AM	COFFEE (10:30-11:00)
11:00 AM	Fundamentals of Counterfactuals in AI
	<i>Formalisation</i>
	<i>Modelling Approaches & Key Constraints</i>
11:30 AM	Using Counterfactual Algorithms
	Hands-on: <i>A Counterfactual Toolbox (AA)</i>
	Hands-on: <i>Checking Out Notebooks and Q&A</i>
12:00 PM	Fundamentals of User Studies
	<i>User Studies I: A Simple Two-Group Design</i>
12:30 PM	LUNCH (12:30-14:00)
2:00 PM	Algorithmic Growth Points
	<i>Computational Future Directions and Q&A</i>
2:30 PM	More Fundamentals of User Studies
	<i>User Studies I: A Simple Two-Group Design (cont.)</i>
3:00 PM	COFFEE (15:00-15:30)
3:30 PM	From Fundamentals to an Actual User Study
	<i>User Studies II: A More Complex Design</i>
	<i>User Studies III: Even More Complex Designs</i>
	Hands-on: <i>Looking At Our Study</i>
5:00 PM	Closing Session, Discussion and Final Q&A
	TUTORIAL END



Other & Current Topics

- **Semi-Factuals**
- Fairness
- Robustness
- Manipulations
- Group Counterfactuals

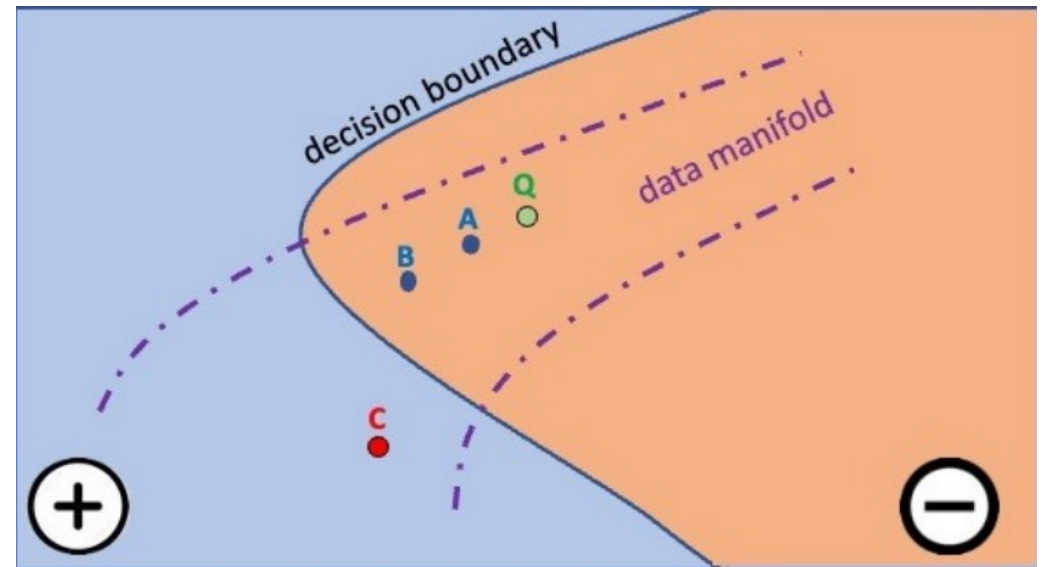


Semi-Factual Explanations

- "Even-If" Explanations [\[Aryal 2023\]](#)
- “**Even if** you used twice as much fertilizer last month, the crop yield **would still have been the same**”
- Why Semi-Factuals?
 - Potentially have major impact (decreasing causal link)
 - May be useful for positive outcomes
 - Have "good" emotional impacts

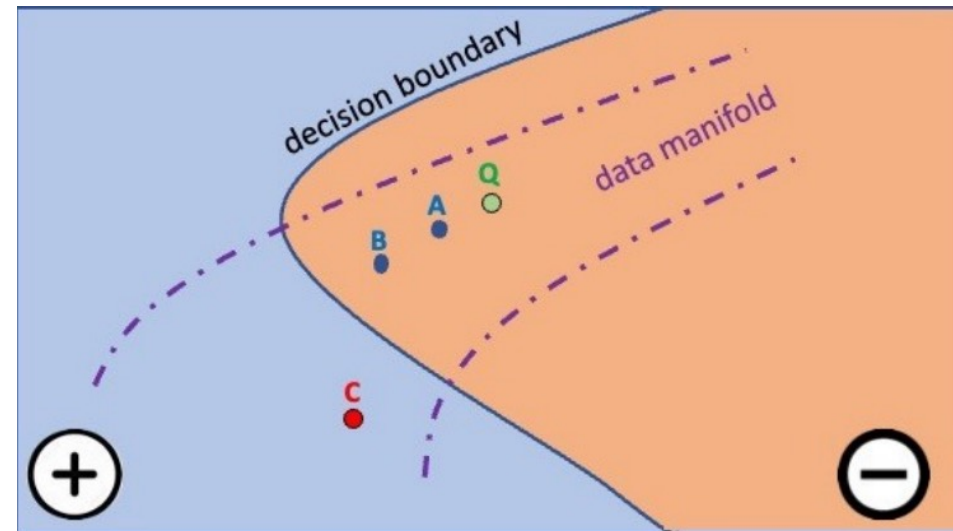
Semi-Factual Explanations

- "Even-If" -- Closer to decision boundary but do not cross it!
- Sub-type of CFs – counter to the facts, still an alternative world



Most Distant Neighbors

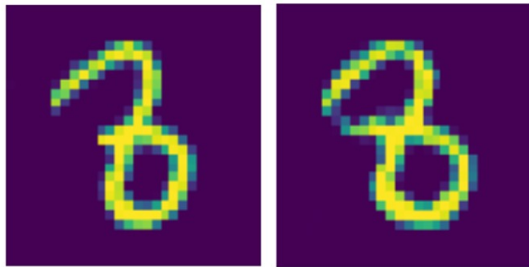
- *Baseline*: Find instance that is furthest from the query on any feature-dimension
- Analogous of NUNs (nb. actual data)
- Does not use counterfactual to guide search (aka CF-Free method)



PIECE

- Plausible Exceptionality-based Contrastive Explanations
- Computes counterfactuals and semi-factuals

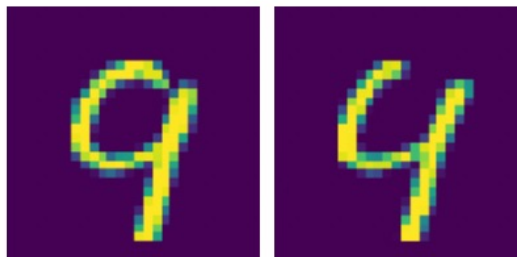
If the test image looked like this, I would have thought it was an "8".



Test Image
Label: 8
Prediction: 3

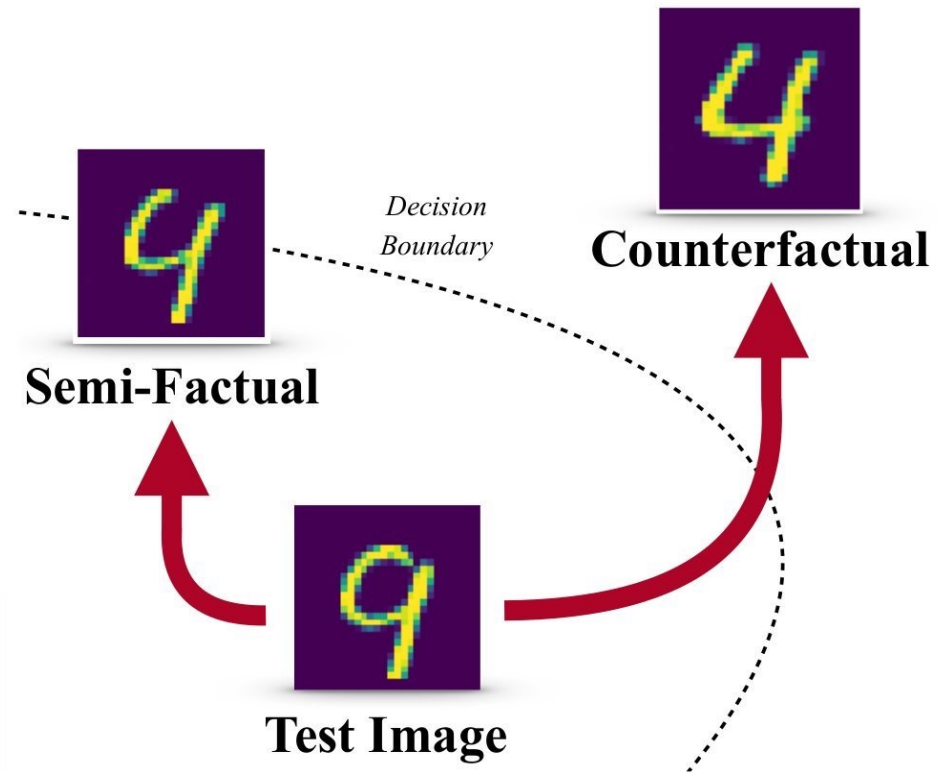
Counterfactual
New Prediction: 8

Even if the test image looked like this, I still would have thought it was a "9".



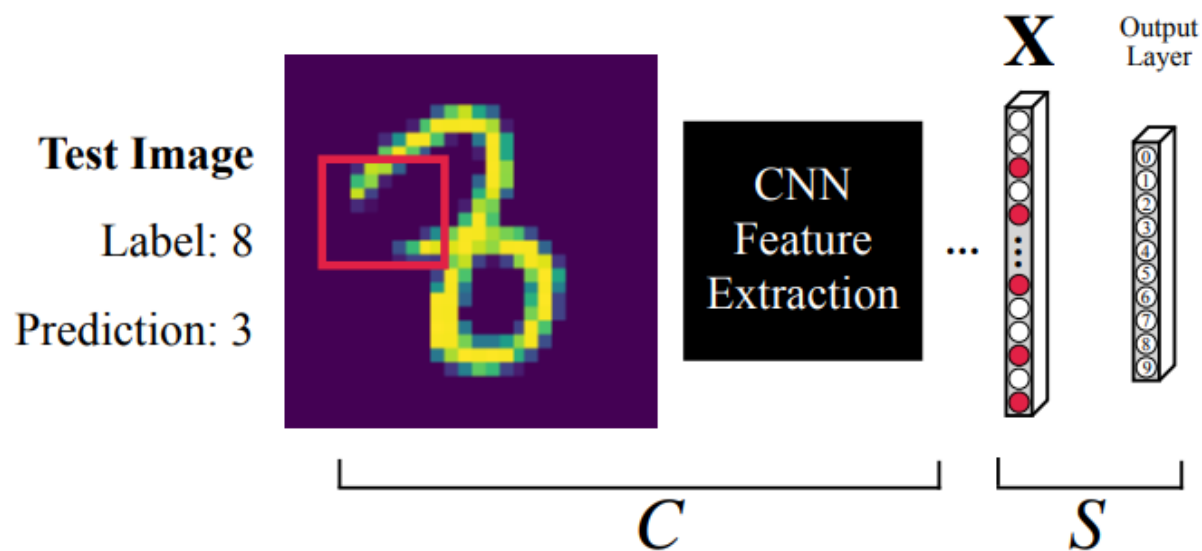
Test Image
Label: 9
Prediction: 9

Semi-Factual
New Prediction: 9

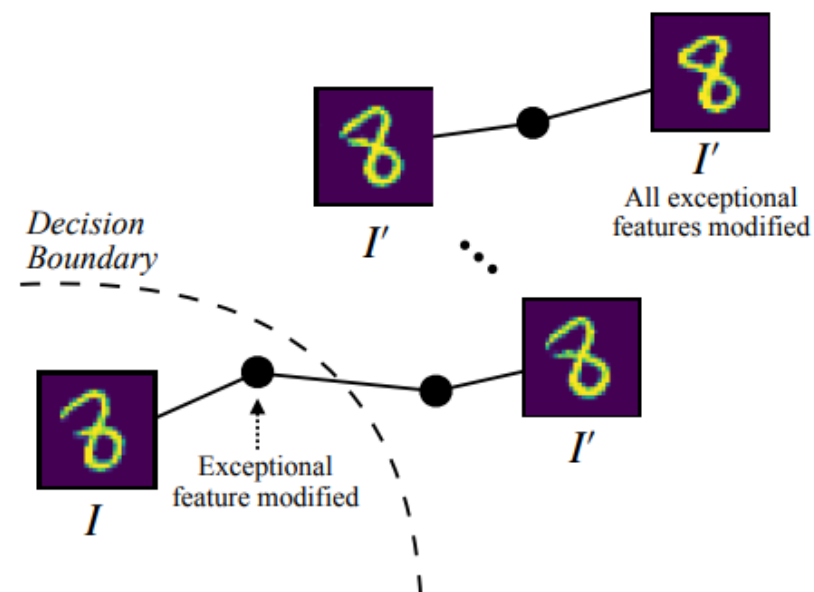


PIECE

(a) Identify Exceptional Features



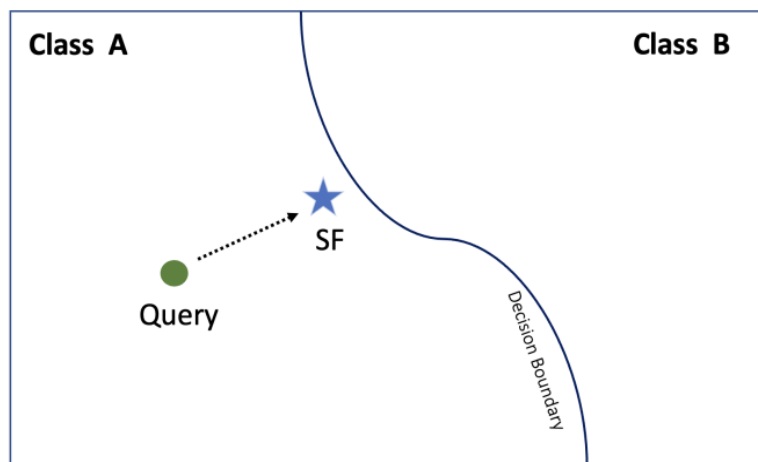
(b) Generate Explanation



Computational Approaches

Counterfactual-Free

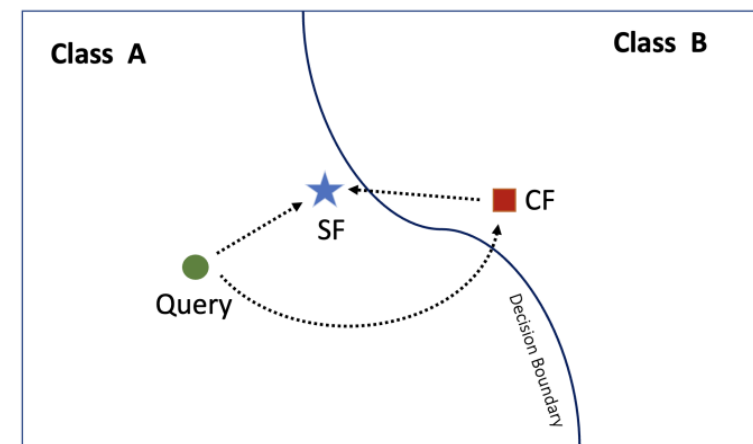
- Diverse Semifactual Explanations of Reject [\[Artelt 2022\]](#)
- Most Distant Neighbors [\[Aryal & Keane 23\]](#)



(a) Counterfactual-Free

Counterfactual-Guided

- PIECE [\[Kenny & Keane AAAI-21\]](#)
- C2C-VAE [\[Zhao et al. ICCBR-22\]](#)



(b) Counterfactual-Guided

Other & Current Topics

- Semi-Factuals
- Fairness
- Robustness
- Manipulations
- Group Counterfactuals



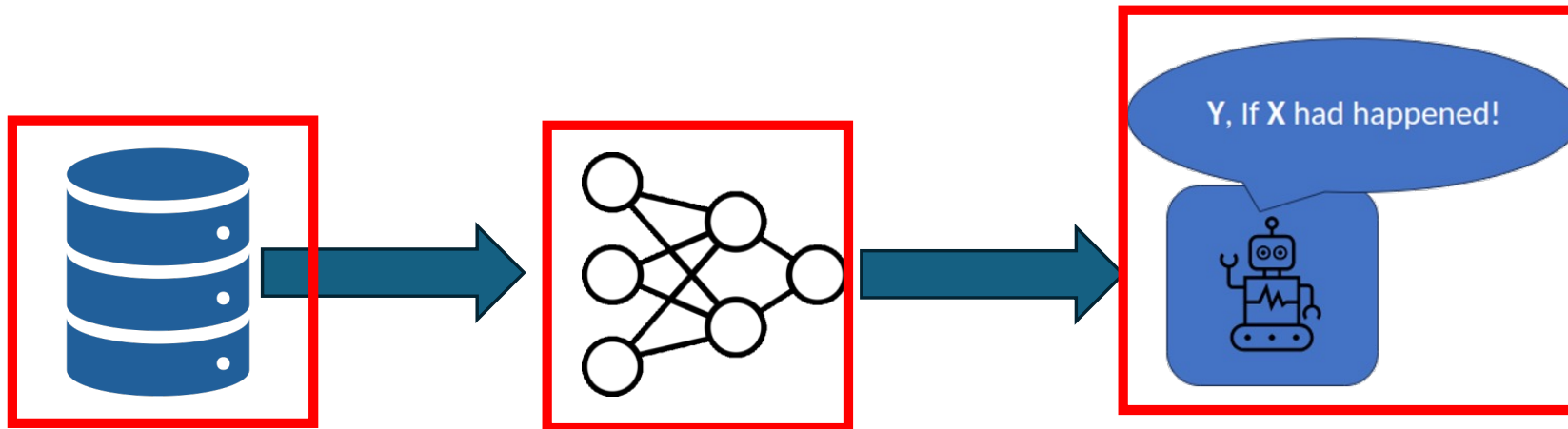
Fairness

- How can a CF be unfair?
=> Differences in the cost of recourse (global vs. local)
- Group vs. Individual fairness [[Slack 2021](#), [Sharma 2020/21](#), [Kügelgen 2022](#), [Artelt 2023](#)]



Approaches to Fix Fairness Issues

- Fix model [[Sharma 2020/21](#)]
- Fix CF generator [[Kügelgen 2022](#), [Artelt 2023](#)]
- Fix data [[Artelt 2024](#)]



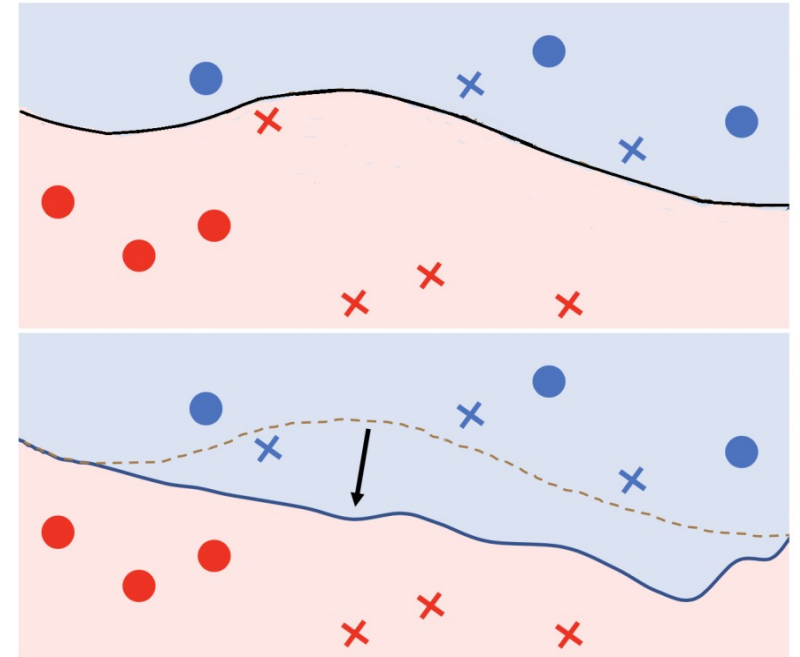
Make the Model Fair [\[Sharma 2020/21\]](#)

- Define cost of recourse for a group "g" as "Burden":

$$Burden(g) = \mathbb{E}_g[d(\mathbf{x}, \mathbf{c}^*)]$$

- Add fairness loss to training loss:

$$\mathcal{L}_{fairness} = \mathbb{E}_{\mathbf{x} | s(\mathbf{x})=a} [d(\mathbf{x}, \mathcal{B})] - \mathbb{E}_{\mathbf{x} | s(\mathbf{x})=b} [d(\mathbf{x}, \mathcal{B})]$$



Fix the CF generator [\[Artelt 2023\]](#)

- Goal: Same distributions of cost of recourse

$$p_{\vec{x}_{orig} \sim \mathcal{D}_1} \left(\theta(\vec{x}_{orig}, CF(\vec{x}_{orig}, h)) \right) \approx p_{\vec{x}_{orig} \sim \mathcal{D}_2} \left(\theta(\vec{x}_{orig}, CF(\vec{x}_{orig}, h)) \right)$$

- Randomized algorithm:

=> Sample cost of recourse

$$\min_{\vec{x}_{cf} \in \mathbb{R}^d} \theta(\vec{x}_{orig}, \vec{x}_{cf}) + C_0 \cdot \ell(h(\vec{x}_{cf}), y_{cf}) + C_1 \cdot \max(z - \theta(\vec{x}_{orig}, \vec{x}_{cf}))$$

where $z \sim p(\theta(\vec{X}_{orig}, \vec{X}_{cf}))$

Fix the CF generator

[Artelt 2023]

- Goal: Same

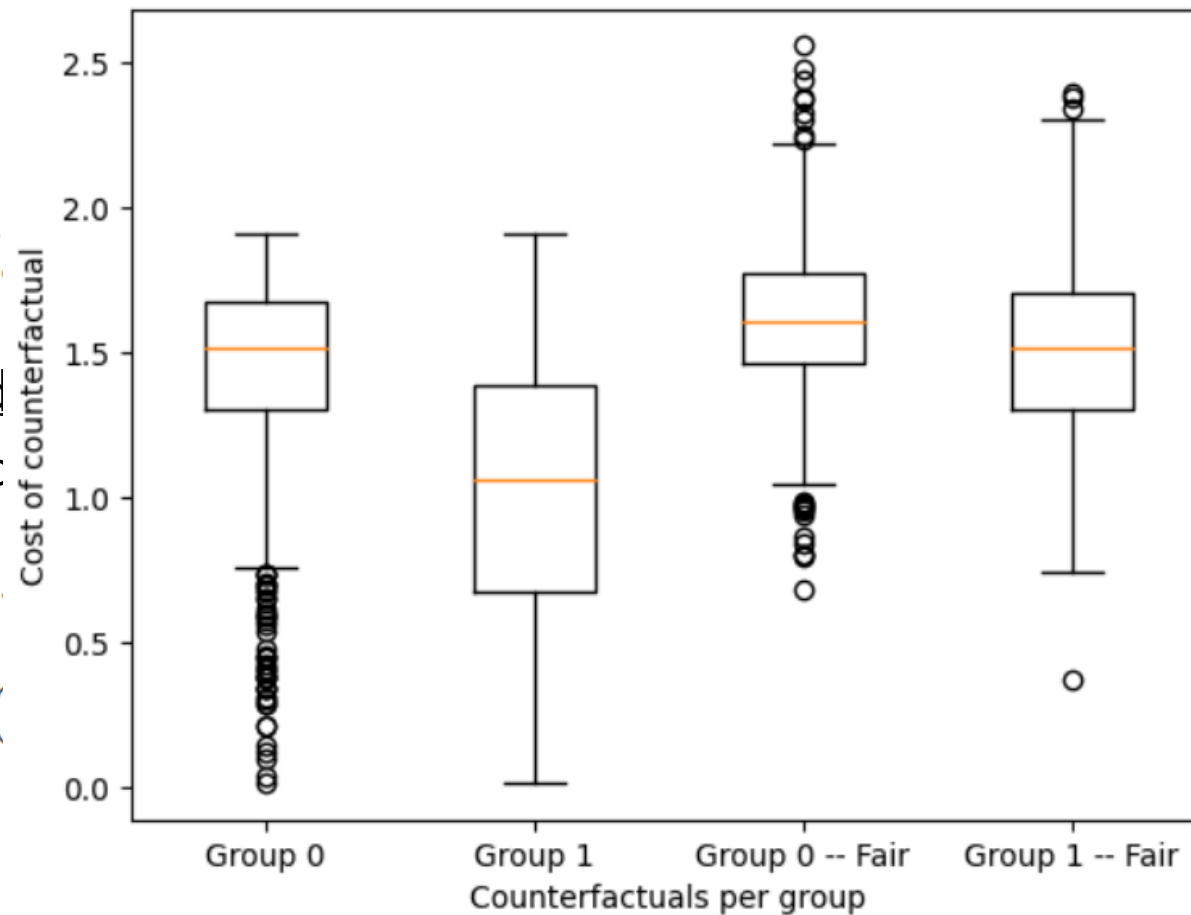
$$p_{\vec{x}_{orig} \sim \mathcal{D}_1} \left(\theta(\vec{x}_{orig}, z) \right)$$

- Randomized

=> Same cc

$$\min_{\vec{x}_{cf} \in \mathbb{R}^d} \theta(\vec{x}_{orig}, z)$$

where $z \sim p(\theta(\vec{x}_{orig}, z))$

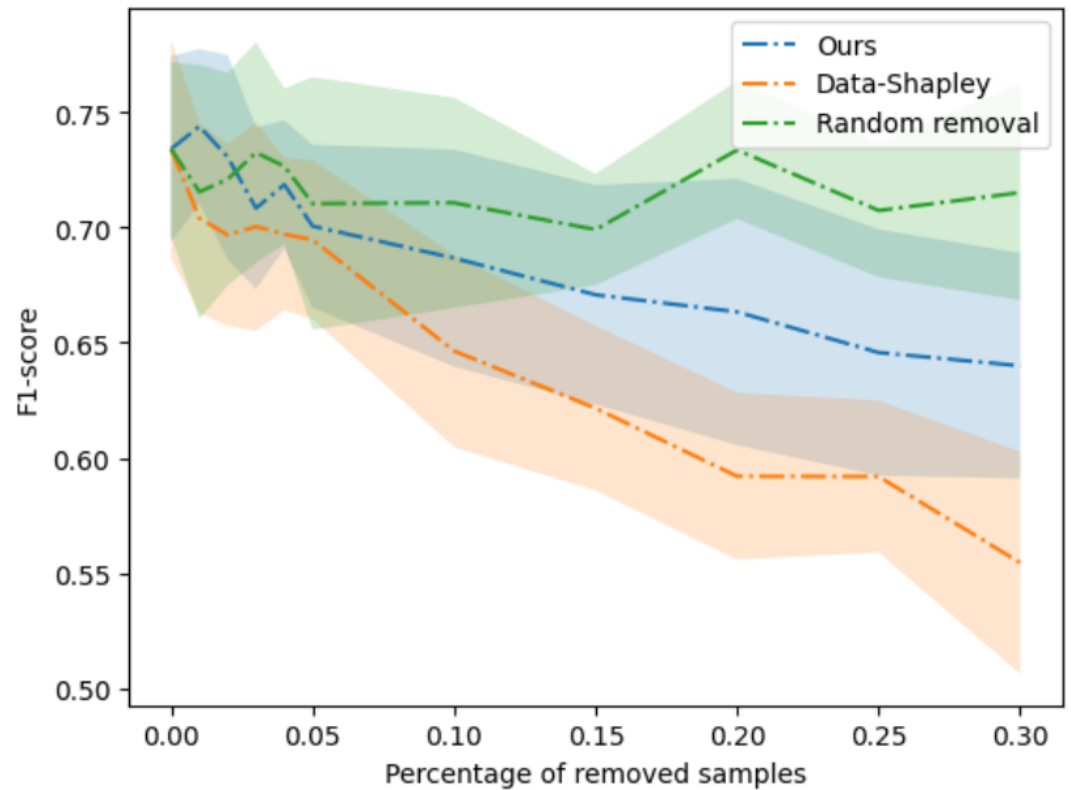
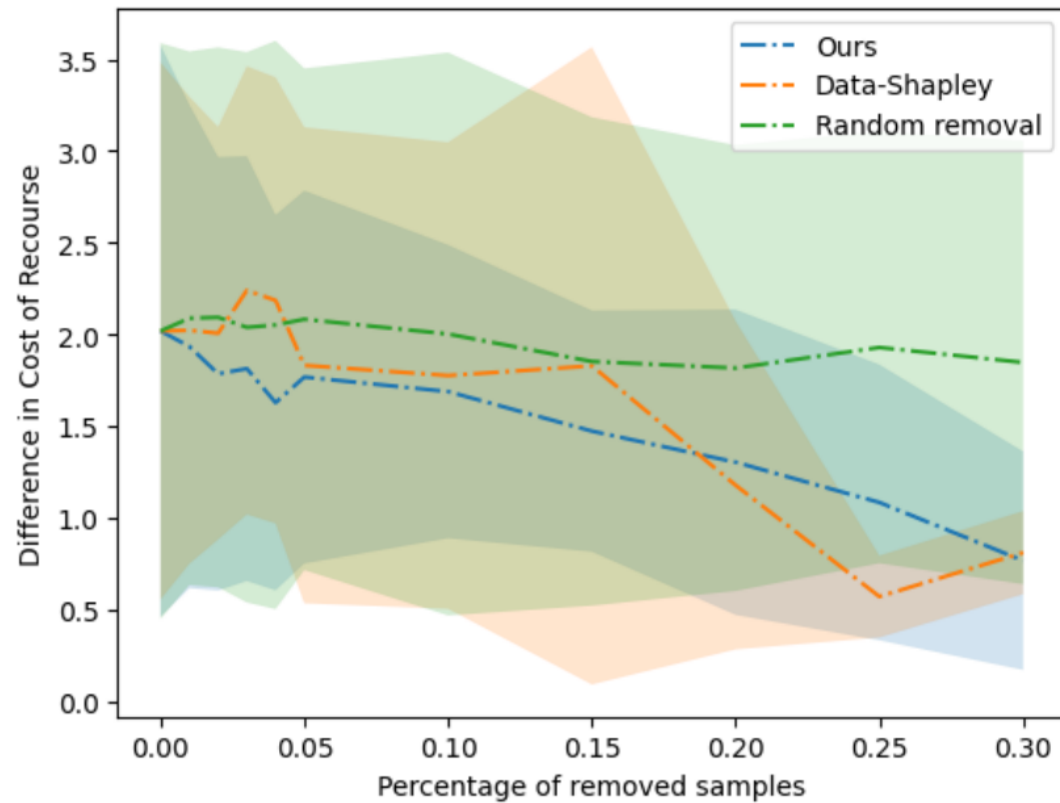


)))

f))

Fix the Data [\[Artelt 2024\]](#)

- Identify influential data points



Open Questions

- Which approach to use?
=> Limitations?
- Formal guarantees?
=> Impossibility results?
- ...

Other & Current Topics

- Semi-Factuals
- Fairness
- Robustness
- Manipulations
- Group Counterfactuals

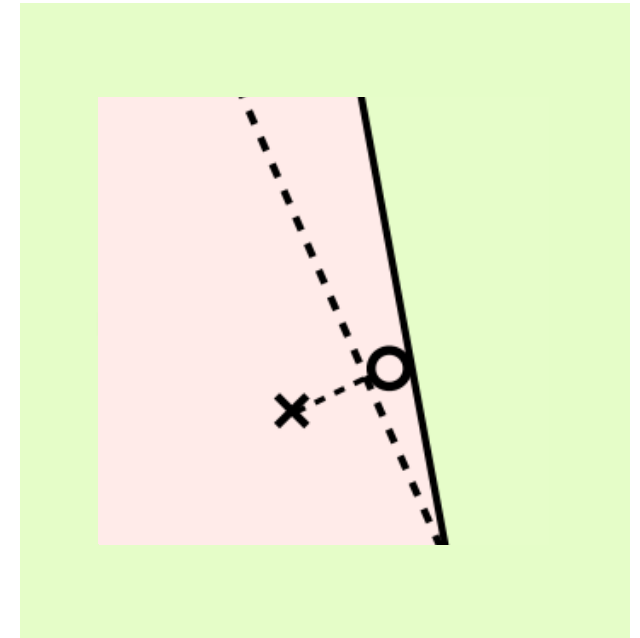


Robustness

- **(Missing) robustness** of CFs [[Jiang 2024](#), [Ferrario 2022](#), [Artelt 2021](#)]
 - Change in explanation
 - Change in cost
- *Input perturbations* (see individual fairness)
- *Model change* (over time)

=> Bad for the user!

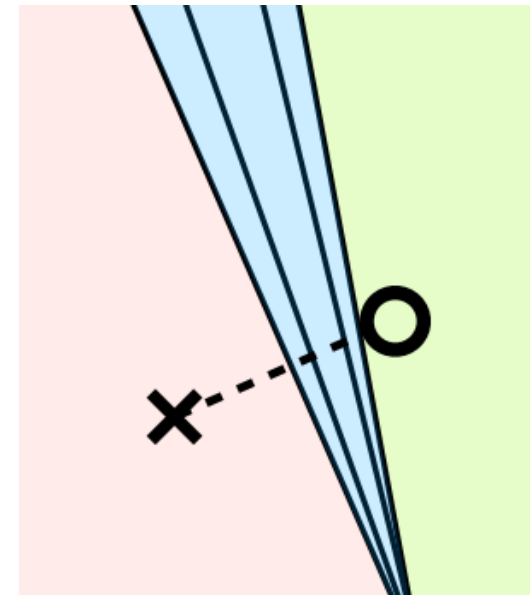
RQ: How to achieve robustness? 🤔



[Leofante 2023]

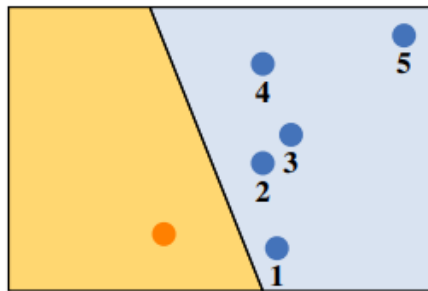
RQ: How to make CFs Robustness? 🤔

- Input changes:
 - Plausibility [[Artelt 2020-21](#), [Zhang 2023](#)]
 - Min-Max Objective [[Dominguez-Olmedo 2022](#)]
 - Diversity [[Leofante 2024](#)]
- Model changes
 - Plausibility [[Pawelczyk 2020](#)]
 - Possible model changes -> MILPs [[Jiang 2023-24](#)]

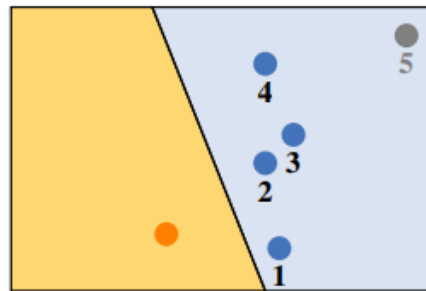


Robustness through Diversity [\[Leofante 2024\]](#)

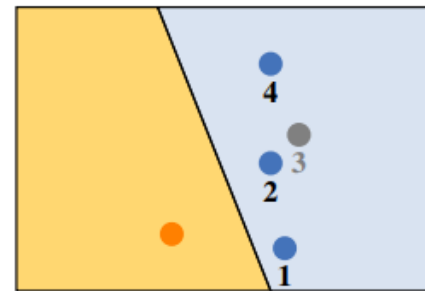
- Compute sets of **diverse CFs** instead of a single one
- **Nearby samples** should have large **overlap** in their diverse CFs!



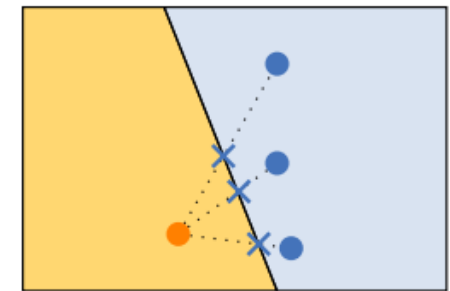
(a)



(b)



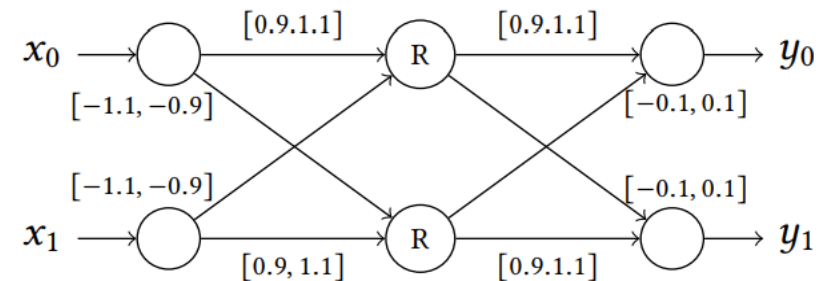
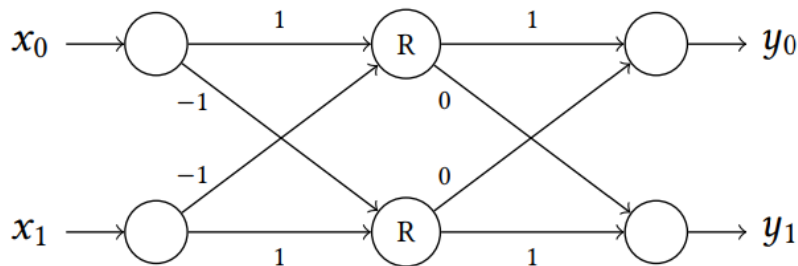
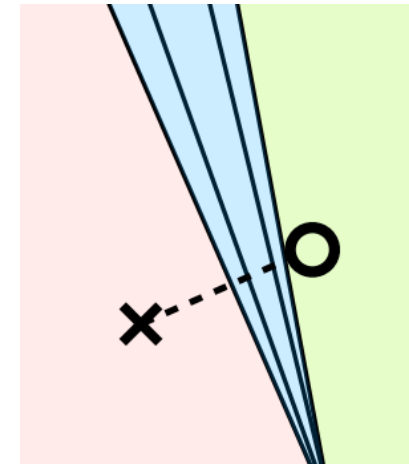
(c)



(d)

Robustness wrt. Model Changes [\[Jiang 2023-24\]](#)

- Consider ranges of parameters
- MILPs for computing & certifying robustness of CFs



Open Questions 🤔

- Cost of robustness
 - How much to sacrifice?
- Robustness and other aspects (e.g. plausibility, fairness, etc.)
 - Impossibility statements?
-

=> Checkout Survey on Robustness @ IJCAI [\[Jiang 2024\]](#)

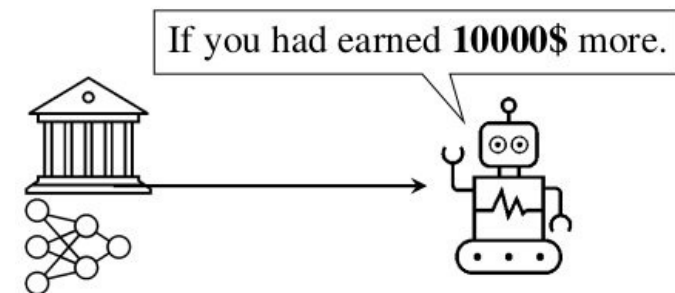
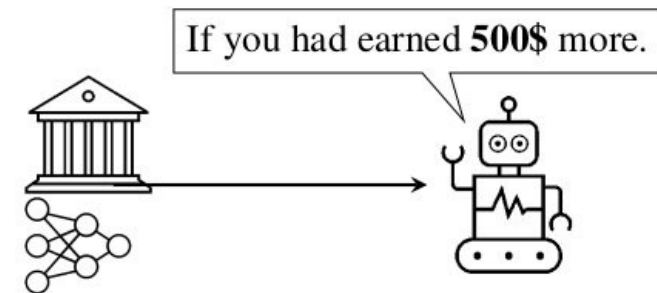
Other & Current Topics

- Semi-Factuals
- Fairness
- Robustness
- Manipulations
- Group Counterfactuals



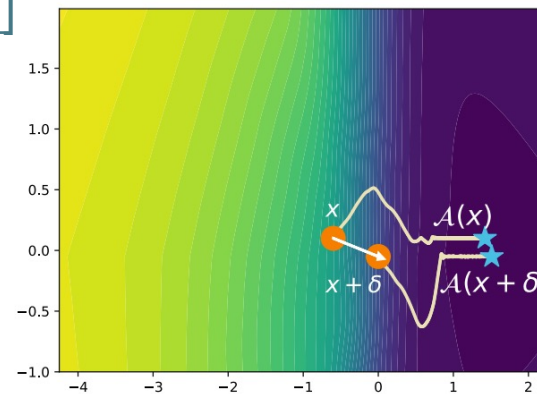
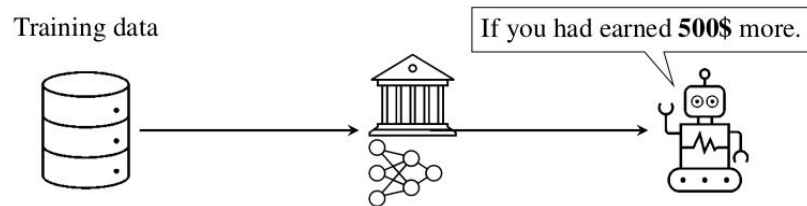
Manipulations

- Can CFs be manipulated?
 - Still trustworthy?
- Attack goals:
 - Cost of recourse
 - Fairness
- Manipulation on which level?

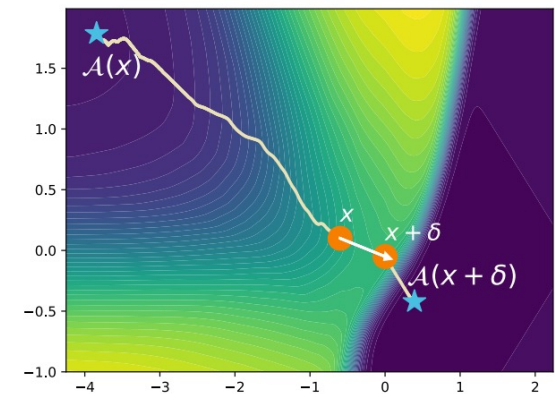


Manipulating CFs

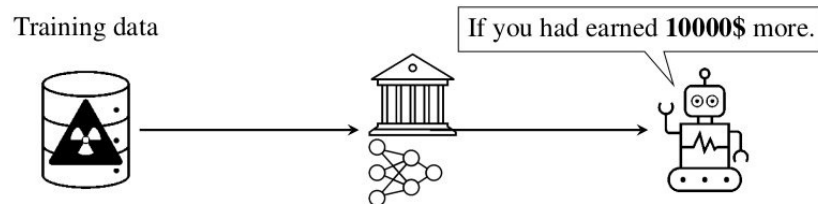
- Adversarial training [\[Slack 2021\]](#)
 - Introduce "unfairness"
- Data poisoning [\[Artelt 2024\]](#)
 - Increase cost of recourse



(a) Training with BCE Objective

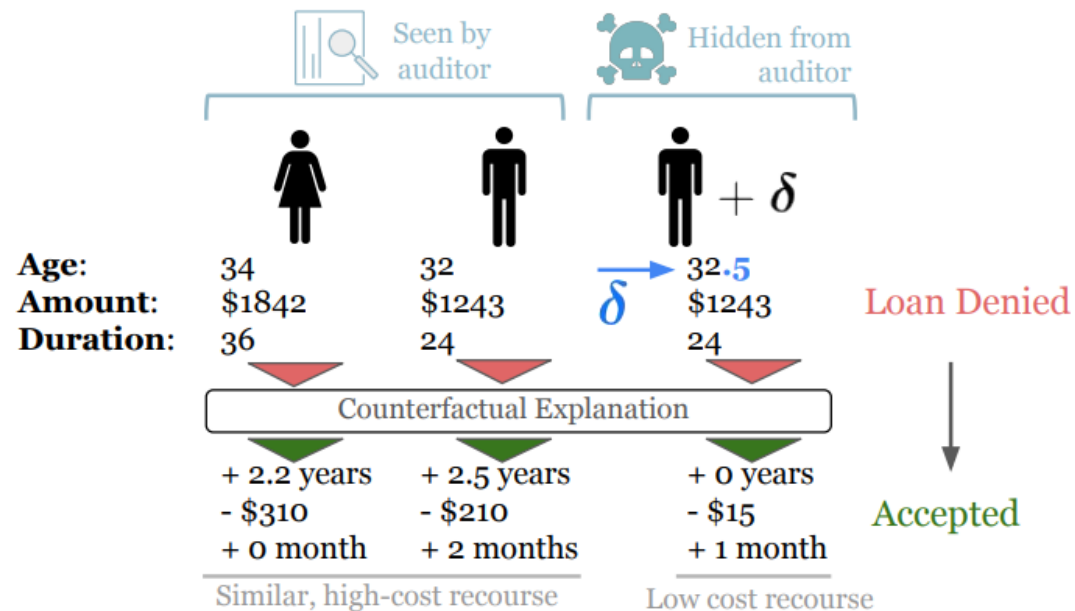


(b) Training Adversarial Model



Adversarial Training [\[Slack 2021\]](#)

- **"Backdoor"**: Some small perturbation leads to lower cost recourse



Adversarial Training [\[Slack 2021\]](#)

- Assume: Attacker can manipulate training process

=> Adversarial training

$$\mathbb{E}_{x \sim \mathcal{D}_{pr}^{neg}} [d(x, \mathcal{A}(x))] \gg \mathbb{E}_{x \sim \mathcal{D}_{np}^{neg}} [d(x, \mathcal{A}(x + \delta))]$$

"Backdoor"

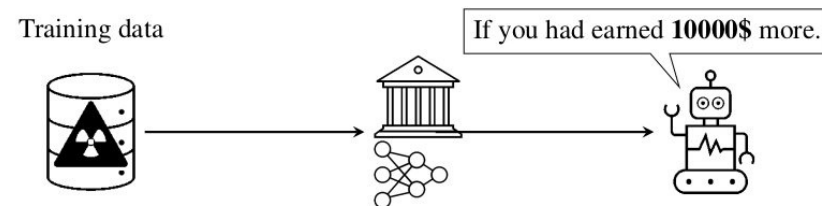
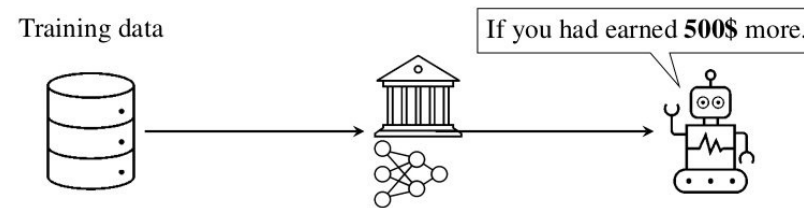
Compute a CF

Table 2: **Recourse Costs of Manipulated Models:** Counterfactual algorithms find similar cost recourses for both subgroups, however, give much lower cost recourse if δ is added before the search.

	Communities and Crime				German Credit			
	Wach.	S-Wach.	Proto.	DiCE	Wach.	S-Wach.	Proto.	DiCE
Protected	35.68	54.16	22.35	49.62	5.65	8.35	10.51	6.31
Non-Protected	35.31	52.05	22.65	42.63	5.08	8.59	13.98	6.81
Disparity	0.37	2.12	0.30	6.99	0.75	0.24	0.06	0.5
Non-Protected+ δ	1.76	22.59	8.50	9.57	3.16	4.12	4.69	3.38
Cost reduction	20.1×	2.3×	2.6×	4.5×	1.8×	2.0×	2.2×	2.0×

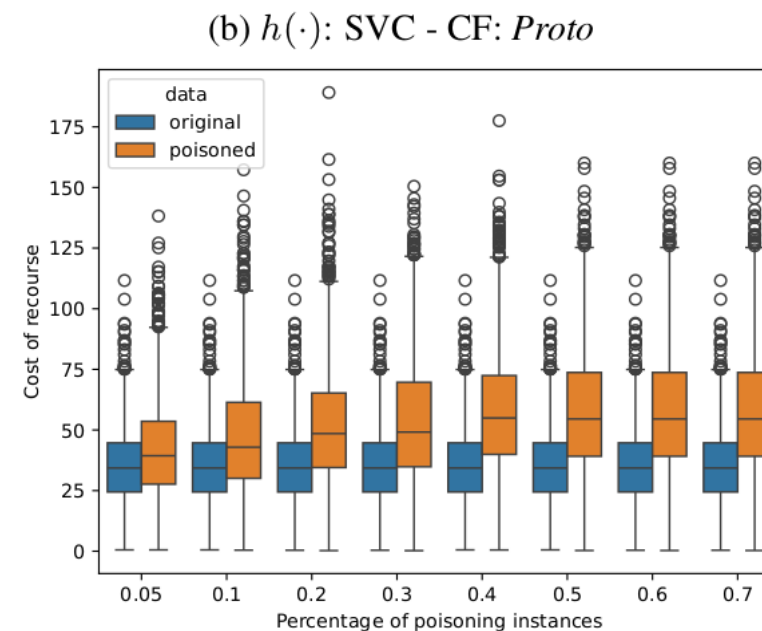
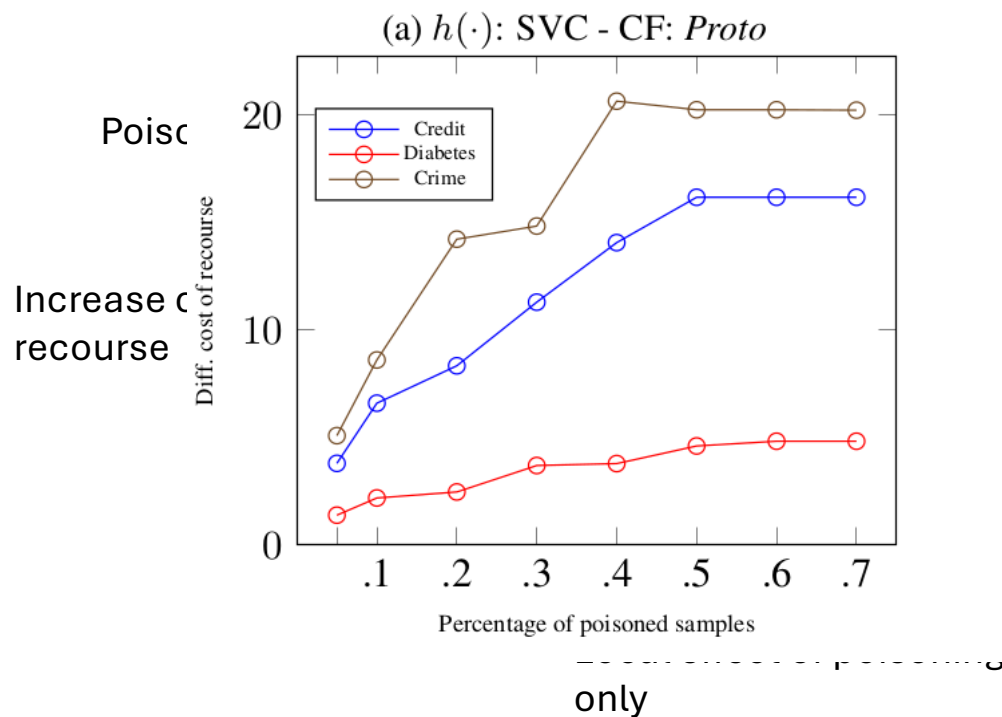
Data Poisoning [\[Artelt 2024\]](#)

- Add poisonous samples to training data
=> Increase cost of recourse



Data Poisoning [\[Artelt 2024\]](#)

- Finding poisonous samples as an optimization problem:



Defense against Manipulations

- Observation: Additional density constraints help [\[Artelt 2024\]](#)
- **Open questions!** 

Other & Current Topics

- Semi-Factuals
- Fairness
- Robustness
- Manipulations
- Group Counterfactuals



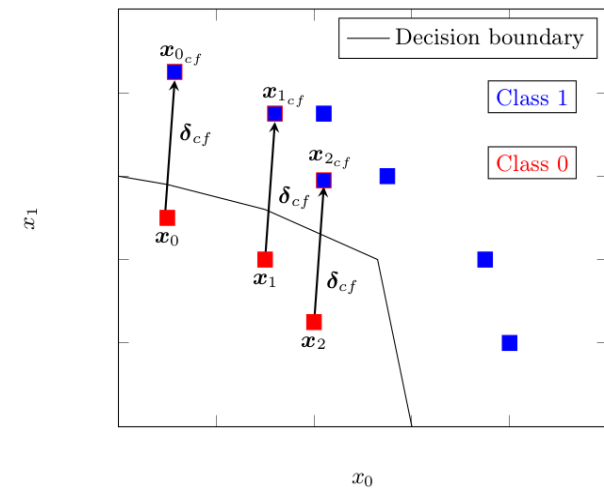
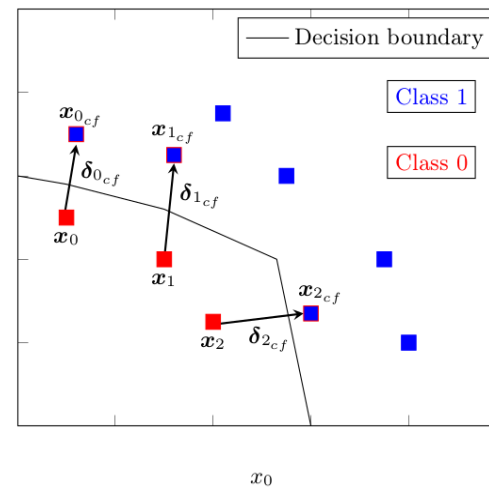
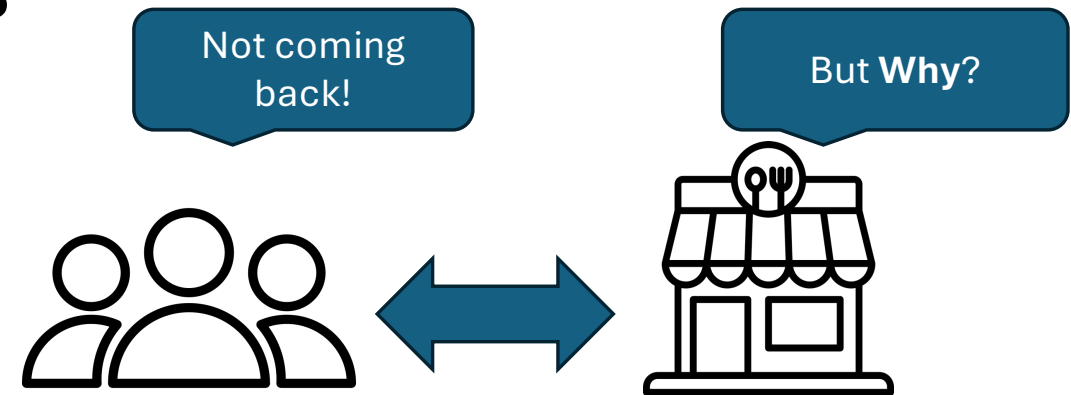
Group Counterfactuals

- One CF for many instances
- Real-world problems, e.g. Customer repurchase [\[Artelt 2024\]](#)

New challenges:

- Clustering/Grouping of instances

Current approaches [\[Kanamori 2022, Ley 2023, Warren 2023, Artelt 2023/24\]](#)



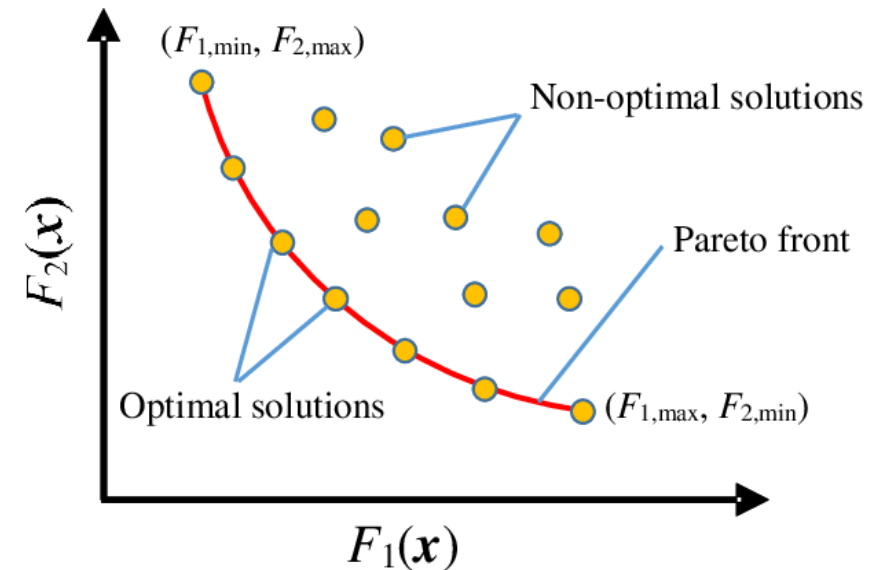
Group Counterfactuals

- Formalization [\[Artelt 2024\]](#):

$$\min_{\vec{\delta}_{\text{cf}} \in \mathbb{R}^d} \left(\underbrace{\theta(\vec{\delta}_{\text{cf}})}_{\text{Cost/Proximity}}, \underbrace{\sum_{\vec{x}_i \in \mathcal{D}} \ell(h(\vec{x}_i \oplus \vec{\delta}_{\text{cf}}), y_{\text{cf}})}_{\text{Contrastive}} \right)$$

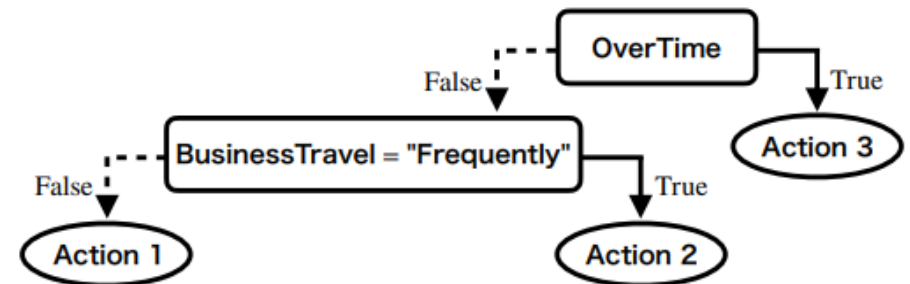
- Feasibility?
=> Pareto-optimal solutions!

What about other constraints? 🤔



Group Counterfactuals

- Grouping given?
- Grouping and CFs
 - All-in-One [\[Kanamori 2022\]](#)
 - Separate [\[Artelt 2024\]](#)



	HowToChange	Effectiveness	
		Cost	Flip rate
Action 1	MonthlyIncome : + 1282\$	0.17	83 %
Action 2	BusinessTravel : "Frequently" → "Rarely"	0.19	80 %
Action 3	OverTime : True → False	0.27	86 %

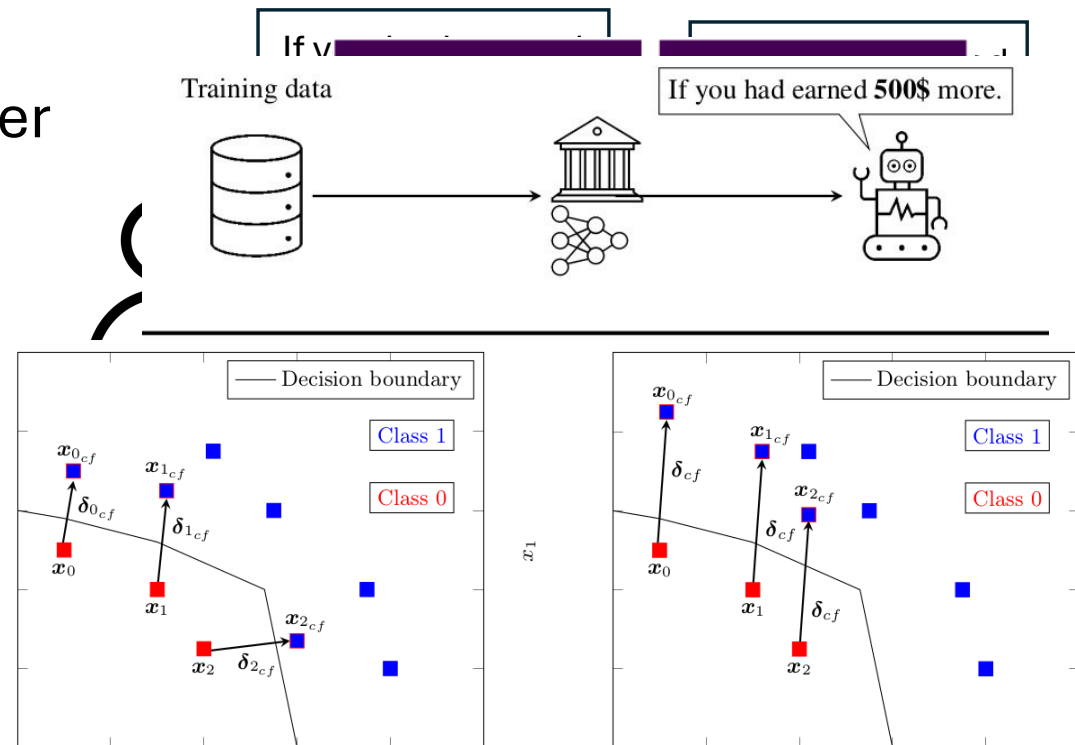
- (a) Cluster with $d(\vec{\delta}_{cf_i}, \vec{\delta}_{cf_j}) = \frac{\vec{\delta}_{cf_i}^\top \vec{\delta}_{cf_j}}{\|\vec{\delta}_{cf_i}\|_2 \|\vec{\delta}_{cf_j}\|_2}$
- (b) Sub-cluster with $d(\vec{\delta}_{cf_i}, \vec{\delta}_{cf_j}) = \|\theta(\vec{\delta}_{cf_i}) - \theta(\vec{\delta}_{cf_j})\|_2$

Open Questions

- Fairness?
- Robustness?
 - => Outliers, poisonous instances, etc.
- Computational complexity and limitations?
 - => Formal analysis
- Python framework/implementation
- ...

Summary

- **Robustness**
- **Fairness**
- **Semi-factuais** as the little brother of CFs
- CFs can be **manipulated!**
 - How to defend?
- CFs for groups of instances -- i.e. **group CFs**
 - Gaining popularity



Questions and Discussions

Ask us anything!